CLAIMS

Having thus described our invention, what we claim as new and desire to secure by Letters Patent is as follows:

1	1. A method of outlier detection comprising the steps of:
2	reducing an outlier detection problem to that of a classification
3	learning problem using unlabeled normal data as positive examples and
4	randomly generated synthesized examples as negative examples, and
5	then selectively sampling normal and synthesized examples based on
6	uncertainty of prediction to further reduce an amount of data required for data
7	analysis, resulting in enhanced predictive performance while minimizing
8	computational resources including storage requirements.
1	2. The method of outlier detection recited in claim 1, wherein the step of
2	selectively sampling includes the step of setting sampling probability equal to
3	a measure of uncertainty of prediction for the example.
1	3. The method of outlier detection recited in claim 2, wherein the measure of
2	uncertainty of prediction is binomial.
1 .	4. The method of outlier detection recited in claim 2, wherein the measure of
2	uncertainty of prediction is Gaussian.
1	5. The method of outlier detection recited in claim 1, wherein the step of
2	selectively sampling includes the step of setting sampling probability
3	proportional to a product of a measure of uncertaint6y and a measure of cust
4	of mis-classifying the same example.

2	6. The method of outlier detection recited in claim 5, wherein the measure of uncertainty is binomial.
1 2	7. The method of outlier detection recited in claim 5, wherein the measure of uncertainty is Gaussian.
1 2 3	8. The method of outlier detection recited in claim 5, wherein the measure of cost is determined as a relative cost of mis-classifying the example in the training data.
1 2 3	9. The method of outlier detection recited in claim 1, wherein the classification learning problem employs an arbitrary algorithm for classification.
1 2 3 4	10. The method of outlier detection recited in claim 9, wherein the arbitrary algorithm for classification is selected from the group consisting of decision tree learning algorithms, naïve Bayes method, logistic regression method and neural network training algorithms.
1 2 3	11. The method of outlier detection recited in claim 1, further comprising the steps of: reading a storing normal data at <i>T</i> -real;
456	generating and storing synthesized data as T -syn; and wherein the step of selective sampling is performed on data $T := T$ -real $\cup T$ -syn.

1	12. The method of outlier detection recited in claim 11, wherein the step of
2	selective sampling uses an underlying, arbitrary classification learning
3	algorithm and proceeds iteratively.
1	13. The method of outlier detection recited in claim 12, wherein each iteration
2	comprises the steps of:
3	selecting a smaller sub-sample from the input data;
4	training of the underlying classification algorithm with the selected
5	data; and
6	storing a classifier output by the classification algorithm.
1	14. The method of outlier detection recited in claim 13, wherein the step of
2	selecting is done by choosing examples that are harder to classify with the
3	classifiers obtained in preceding iterations.
1	15. The method of outlier detection recited in claim 13, further comprising the
2	step of outputting an output hypothesis as a voting function of classifiers
3	obtained in the iterations.
1	16 The mosths of a Carellin 1 and the state of the state
	16. The method of outlier detection recited in claim 13, wherein the step of
2	selecting is done by choosing each example with a sampling probability which
3	is set equal to a measure of uncertainty of predicting a label of that example by
4	a collection of hypotheses obtained by calls to the classification algorithm in
5	earlier iterations.
1	17. The method of outlier detection recited in claim 13, wherein the step of
2	selecting is done by choosing each example with a sampling probability which
3	is set proportional to a product of a measure of uncertainty of predicting a

4	later of that example by a collection of hypotheses obtained by calls to the
5	classification algorithm in earlier iterations and a measure of cost of mis-
6	classifying the same example.
1	18. A data processing system for outlier detection comprising:
2	a top control module controlling overall control flow, making use of
3	various sub-components of the system;
4	a learning algorithm storage module storing a representation of an
5	algorithm for classification learning;
. 6	a model output module storing models obtained as a result of applying
7	the learning algorithm stored in learning algorithm storage module to training
8	data and outputting a final model by aggregating these models; and
9	a selective sampling module accessing data stored in a data storage
10	module, selectively sampling a relatively small subset of the data, and passing
11	the obtained sub-sample to the top control module.
	10.77
1	19. The data processing system for outlier detection recited in claim 18,
2	wherein the learning algorithm storage module stores an arbitrary algorithm
3	for classification.
1	20. The data processing system for outlier detection recited in claim 19,
2	wherein the arbitrary algorithm for classification is selected from the group
3 —	consisting of decision tree learning algorithms, naïve Bayes method, logistic
4	regression method and neural network training algorithms.
1	21. The data processing system for outlier detection recited in claim 18,
2	wherein the data storage module comprises two separate modules, one for

- storing real data corresponding to "normal" data, and the other for storing
 synthesized data corresponding to "abnormal" data.
- 1 22. The data processing system for outlier detection recited in claim 18,
- 2 wherein the data storage module comprises a single data storage module
- 3 providing two logical data storage modules in a single physical data storage
- 4 module, one for storing real data corresponding to "normal" data, and the
- 5 other for storing synthesized data corresponding to "abnormal" data.